

2012

# Privacy-Preserving Data Sharing in High Dimensional Regression and Classification Settings

Stephen E. FIENBERG

Jiashun JIN

Follow this and additional works at: <https://ink.library.smu.edu.sg/larc>

Part of the [Databases and Information Systems Commons](#), and the [Information Security Commons](#)

---

## Citation

FIENBERG, Stephen E. and JIN, Jiashun. Privacy-Preserving Data Sharing in High Dimensional Regression and Classification Settings. (2012). LARC Research Publications.

**Available at:** <https://ink.library.smu.edu.sg/larc/1>

This Report is brought to you for free and open access by the School of Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in LARC Research Publications by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email [libIR@smu.edu.sg](mailto:libIR@smu.edu.sg).



**Privacy-Preserving Data Sharing in High Dimensional Regression and Classification Settings**

***Stephen E. Fienberg, Carnegie Mellon University, Pittsburgh, PA USA***

*fienberg@stat.cmu.edu*

***Jiashun Jin, Carnegie Mellon University, Pittsburgh, PA USA***

*jiashun@stat.cmu.edu*

August, 2012

LARC-TR-05-12

LARC Technical Report Series: <http://smu.edu.sg/centres/larc/larc-technical-reports-series/>



## ***ABSTRACT***

We focus on the problem of multi-party data sharing in high dimensional data settings where the number of measured features (or the dimension)  $p$  is frequently much larger than the number of subjects (or the sample size)  $n$ , the so-called  $p \gg n$  scenario that has been the focus of much recent statistical research. Here, we consider data sharing for two interconnected problems in high dimensional data analysis, namely the feature selection and classification. We characterize the notions of "cautious", "regular", and "generous" data sharing in terms of their privacy-preserving implications for the parties and their share of data, with focus on the "feature privacy" rather than the "sample privacy," though the violation of the former may lead to the latter. We evaluate the data sharing methods using a *phase diagram* from the statistical literature on multiplicity and Higher Criticism thresholding. In the two-dimensional phase space calibrated by the signal sparsity and signal strength, a phase diagram is a partition of the phase space and contains three distinguished regions, where we have no (feature) privacy violation, relatively rare privacy violations, and an overwhelming amount of privacy violation.

This report has been published in the Journal of Privacy and Confidentiality. The full text of the article can be found here:  
<http://repository.cmu.edu/jpc/vol4/iss1/10>